

COMP307/AIML420

INTRODUCTION TO

ARTIFICIAL INTELLIGENCE



Reasoning under uncertainty:
Bayesian Networks

Outline

- Review of Bayes Theorem and Naïve Bayes
- Introduction to Bayesian Networks

Review

Product rule: $P(A, B) = P(B) * P(A | B) = P(A) * P(B | A)$

Bayes theorem:

- Provides a way to calculate the probability of a hypothesis (e.g., **label**) given some evidence (e.g., **feature values**).

$$P(A | B) = \frac{P(B | A) * P(A)}{P(B)}$$

Naïve Bayes:

- Probabilistic classifier
- Training: count and store priors and likelihoods
- Assumes features are **conditionally independent given the class label**

Review

Product rule: $P(A, B) = P(B) * P(A | B) = P(A) * P(B | A)$

Bayes theorem:

- Provides a way to calculate the probability of a hypothesis (e.g., **class label**) given some evidence (e.g., **feature values**).

$$P(A | B) = \frac{P(B | A) * P(A)}{P(B)}$$

The diagram illustrates the relationship between the terms in Bayes' theorem and the concepts used in Naïve Bayes. A red line connects the prior probability $P(A)$ to the term 'priors' in the Naïve Bayes list. A green line connects the likelihood $P(B|A)$ to the term 'likelihoods' in the Naïve Bayes list. A red line also connects the denominator $P(B)$ to the 'likelihoods' term, indicating that the marginal probability of the evidence is also related to the likelihoods in this context.

Naïve Bayes:


- Probabilistic classifier
- Training: count and store priors and likelihoods
- Assumes features are **conditionally independent given the class label**

Review

- Naïve Bayes:

$$p(c | x_1, \dots, x_p) = \frac{p(c)p(x_1, \dots, x_p | c)}{p(x_1, \dots, x_p)}$$

class conditional feature independence


$$= \frac{p(c)p(x_1 | c) \cdots p(x_p | c)}{p(x_1, \dots, x_p)}$$

- We don't need denominator to find c with highest probability
 - Use *score*

Review: Naïve Bayes

- In classification we know the **priors** and the **likelihoods** for the training data
- Example:
 - Class $C = \text{sunny}$ (true/false)
 - Features: season (spring/summer/fall/winter), humidity (high/low)
 - Want to compute determine
if sunny=true or sunny=false given that *humidity = high, season = summer*
 - Probabilities needed:
 - $P(\text{season}=\text{summer} \mid \text{sunny}=\text{true})$: count instances where season = summer and class label sunny = true and divide by total number of instances where sunny = true
 - $P(\text{season}=\text{summer} \mid \text{sunny}=\text{false})$
 - $P(\text{humidity}=\text{high} \mid \text{sunny}=\text{true})$
 - $P(\text{humidity}=\text{high} \mid \text{sunny}=\text{false})$
 - $P(\text{sunny}=\text{true})$: count the instances sunny=true and divide by the total number of instances
 - $P(\text{sunny}=\text{false})$
- What is the conditional independence assumption? Is it reasonable?
- Why do we calculate a *score* instead of the posterior probability $P(A|B)$?

Review: Naïve Bayes

Given an instance $x = (\mathbf{season=Summer, humidity=high})$, we need to calculate:

$$P(\mathbf{sunny} = \mathit{True} \mid \mathbf{season} = \mathit{Summer}, \mathbf{humidity} = \mathit{high})$$

$$P(\mathbf{sunny} = \mathit{False} \mid \mathbf{season} = \mathit{Summer}, \mathbf{humidity} = \mathit{high})$$

Review: Naïve Bayes

Given an instance $x = (\text{season}=\text{Summer}, \text{humidity}=\text{high})$, we need to calculate:

$$P(\text{sunny} = \text{True} \mid \text{season} = \text{Summer}, \text{humidity} = \text{high})$$

$$P(\text{sunny} = \text{False} \mid \text{season} = \text{Summer}, \text{humidity} = \text{high})$$

*

$$\begin{aligned} P(\text{sunny} = \text{True} \mid \text{season} = \text{Summer}, \text{humidity} = \text{high}) = \\ P(\text{season} = \text{Summer}, \text{humidity} = \text{high} \mid \text{sunny} = \text{True}) * \\ P(\text{sunny} = \text{True}) / \\ P(\text{season} = \text{Summer}, \text{humidity} = \text{high}) \end{aligned}$$

Review: Naïve Bayes

Given an instance $x = (\text{season}=\text{Summer}, \text{humidity}=\text{high})$, we need to calculate:

$$P(\text{sunny} = \text{True} \mid \text{season} = \text{Summer}, \text{humidity} = \text{high})$$

$$P(\text{is_sunny} = \text{False} \mid \text{season} = \text{Summer}, \text{humidity} = \text{high})$$

$$P(\text{sunny} = \text{True} \mid \text{season} = \text{Summer}, \text{humidity} = \text{high}) =$$

$$P(\text{season} = \text{Summer}, \text{humidity} = \text{high} \mid \text{sunny} = \text{True}) *$$

$$P(\text{sunny} = \text{True}) /$$

$$P(\text{season} = \text{Summer}, \text{humidity} = \text{high})$$

$$P(\text{season} = \text{Summer} \mid \text{sunny} = \text{True}) P(\text{humidity} = \text{high} \mid \text{sunny} = \text{True})$$

Conditional independence assumption. Is it reasonable?

Review: Naïve Bayes

Given an instance $x = (\text{season}=\text{Summer}, \text{humidity}=\text{high})$, we need to calculate:

$$P(\text{sunny} = \text{True} \mid \text{season} = \text{Summer}, \text{humidity} = \text{high})$$

$$P(\text{is_sunny} = \text{False} \mid \text{season} = \text{Summer}, \text{humidity} = \text{high})$$

$$P(\text{sunny} = \text{True} \mid \text{season} = \text{Summer}, \text{humidity} = \text{high}) =$$

$$P(\text{season} = \text{Summer}, \text{humidity} = \text{high} \mid \text{sunny} = \text{True}) *$$

$$P(\text{sunny} = \text{True}) /$$

$$P(\text{season} = \text{Summer}, \text{humidity} = \text{high})$$

calculating $P(X)$ not needed for class label

$$P(\text{season} = \text{Summer} \mid \text{sunny} = \text{True}) P(\text{humidity} = \text{high} \mid \text{sunny} = \text{True})$$

conditional independence assumption

Review: Naïve Bayes

Given an instance $x = (\text{season}=\text{Summer}, \text{humidity}=\text{high})$, we need to calculate:

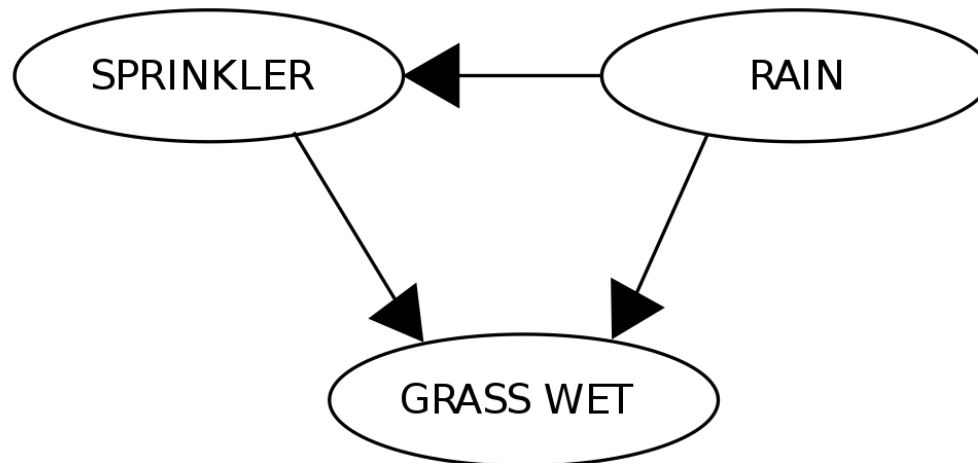
$$P(\text{sunny} = \text{true} \mid \text{season} = \text{summer}, \text{humidity} = \text{high})$$

$$P(\text{sunny} = \text{false} \mid \text{season} = \text{summer}, \text{humidity} = \text{high})$$

$$\begin{aligned} \text{score}(\text{sunny} = \text{true} \mid \text{season} = \text{summ}, \text{hum} = \text{high}) = \\ P(\text{season} = \text{summ} \mid \text{sunny} = \text{true}) * \\ P(\text{hum} = \text{high} \mid \text{sunny} = \text{true}) * \\ P(\text{sunny} = \text{true}) \end{aligned}$$

Bayesian Networks

- **Bayesian networks (BNs)** are a type of **probabilistic graphical model** that represents the joint probability distribution over a set of random variables and their **conditional dependencies** using a directed acyclic graph (DAG).
- **Node:** a random variable
- **Edge:** represent *causal* dependencies between nodes



BN Example 1

Given an electric fan, suppose you try to turn it on, but it doesn't spin (not working).

Why is the fan not spinning?

- **Faulty fan:** the fan is broken
- **Faulty plug:** the plug is broken
- A **phone charger** connected to the same plug works well
- ***“Faulty Fan and Faulty Plug are marginally independent; however, they become conditionally dependent, given Fan.” [1]***

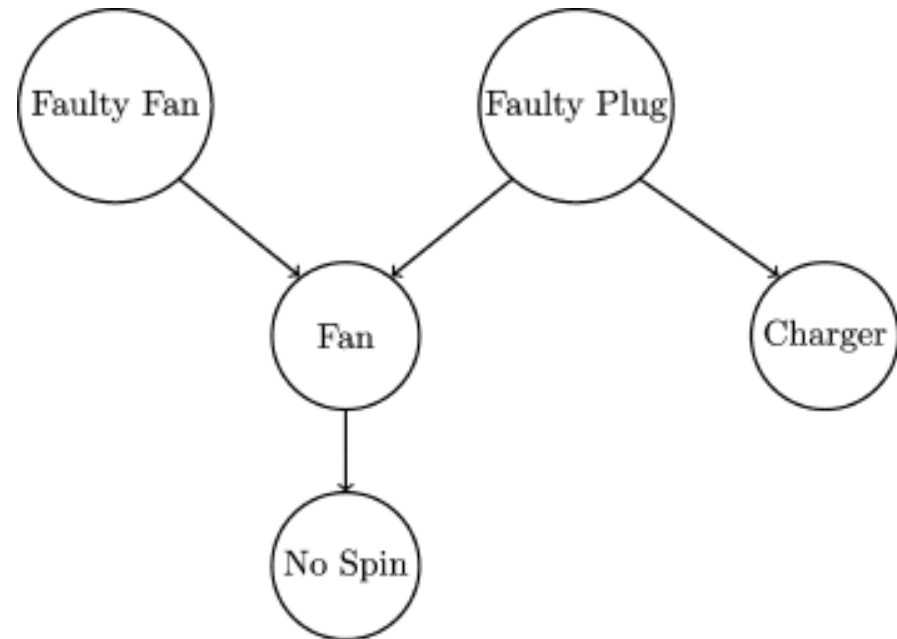


Figure: Simple BN [1]

BN Example 1

Given an electric fan, suppose you try to turn it on, but it doesn't spin (not working).

Why is the fan not spinning?

- **Faulty fan:** the fan is broken
- **Faulty plug:** the plug is broken
- A **phone charger** connected to the same plug works well
- *“**Faulty Fan** and **Faulty Plug** are marginally independent; however, they become conditionally dependent, given Fan.” [1]*
- Important concepts:
 - **(marginally) independent:** if we know the (marginal) probability distribution of one variable, it does not affect the probability distribution of the other variable
 - **conditionally independent:** the variables are independent when another variable is known

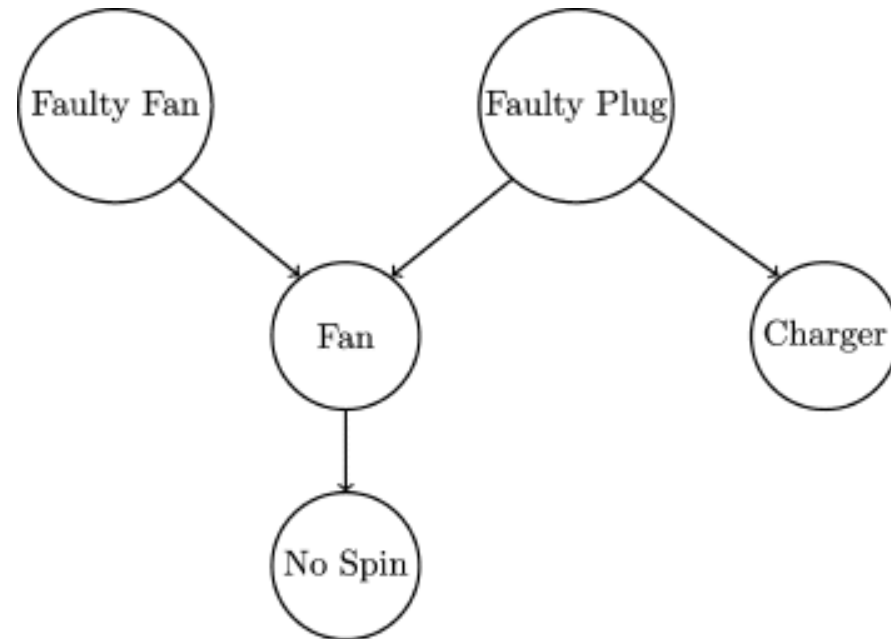


Figure: Simple BN [1]

BN Example 1

Given an electric fan, suppose you try to turn it on, but it doesn't spin (not working).

Why is the fan not spinning?

- **Faulty fan:** the fan is broken
- **Faulty plug:** the plug is broken
- A **phone charger** connected to the same plug works well
- *“**Faulty Fan** and **Faulty Plug** are marginally independent; however, they become conditionally dependent, given Fan.” [1]*
- Important concepts:
 - **(marginally) independent:** if we know the (marginal) probability distribution of one variable, it does not affect the probability distribution of the other variable
 - **conditionally independent:** the variables are independent when another variable is known

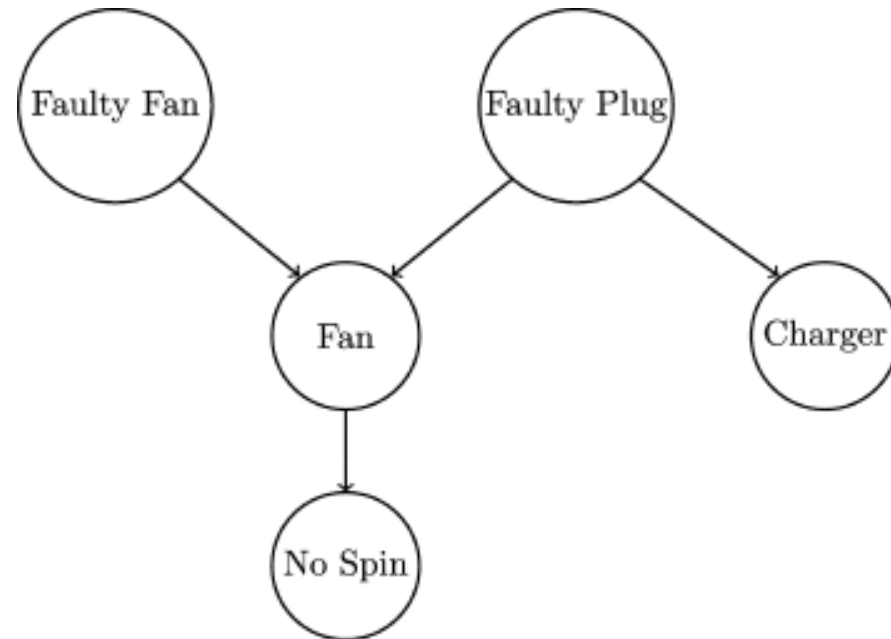


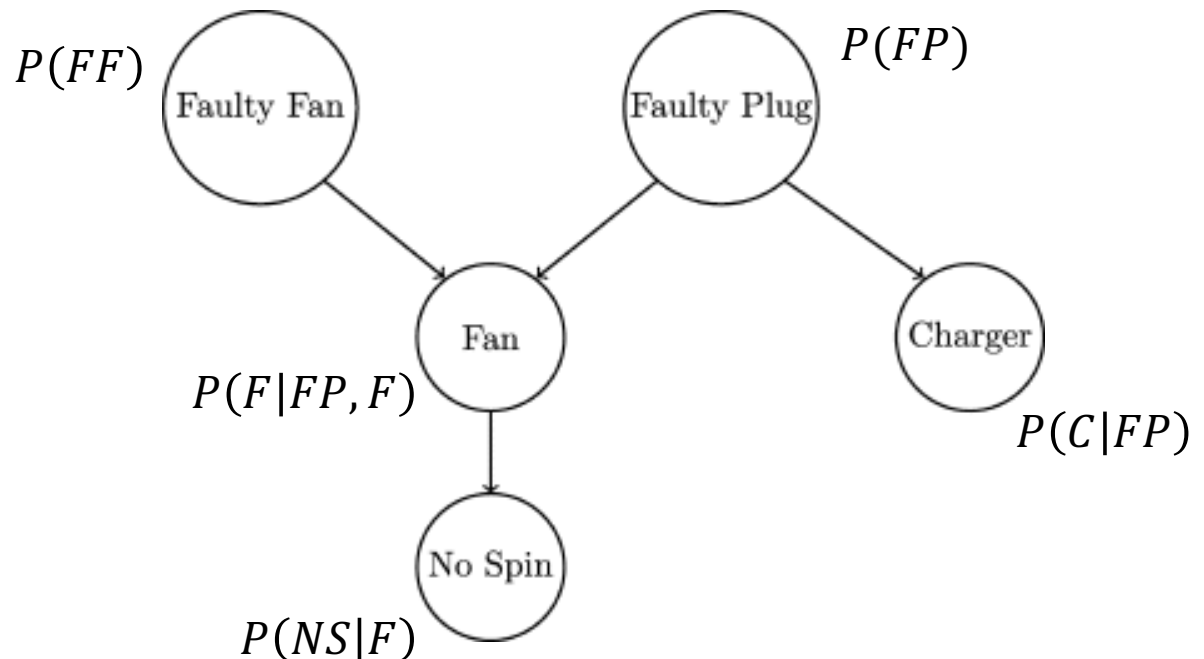
Figure: Simple BN [1]

BN models the joint probability distribution of a set of random variables by decomposing it into a product of conditional probabilities

BN Example 1

- Work backward from the bottom to obtain:

$$\begin{aligned} P(FF, FP, F, C, NS) &= P(NS | FF, FP, F, C) P(FF, FP, F, C) \\ &= P(NS | F) P(F | FP, F, C) P(FF, FP, C) \\ &= P(NS | F) P(F | FP, F) P(C | FF, FP) P(FF, FP) \\ &= P(NS | F) P(F | FP, F) P(C | FP) P(FF) P(FP) \end{aligned}$$



BN Example

$$\begin{aligned}P(FF, FP, F, C, NS) &= P(NS | FF, FP, F, C) P(FF, FP, F, C) \\ &= P(NS | F) P(F | FP, F, C) P(FF, FP, C) \\ &= P(NS | F) P(F | FP, F) P(C | FF, FP) P(FF, FP) \\ &= P(NS | F) P(F | FP, F) P(C | FP) P(FF) P(FP)\end{aligned}$$

- All features have two states. Five features. Joint probability $P(FF, FP, F, C, NS)$ has table with 32 entries (31 sufficient).
- Instead, we have five tables with $2+4+2+1+1 = 10$ entries.
 - For example, $P(NS|F)$ has 2 entries; exploit that $P(NS = 0|F) = 1 - P(NS = 1|F)$ so store only one of these.
- Fewer parameters, hence fewer data needed for their estimation. Parameters replaced with structural knowledge.
 - More robust results.

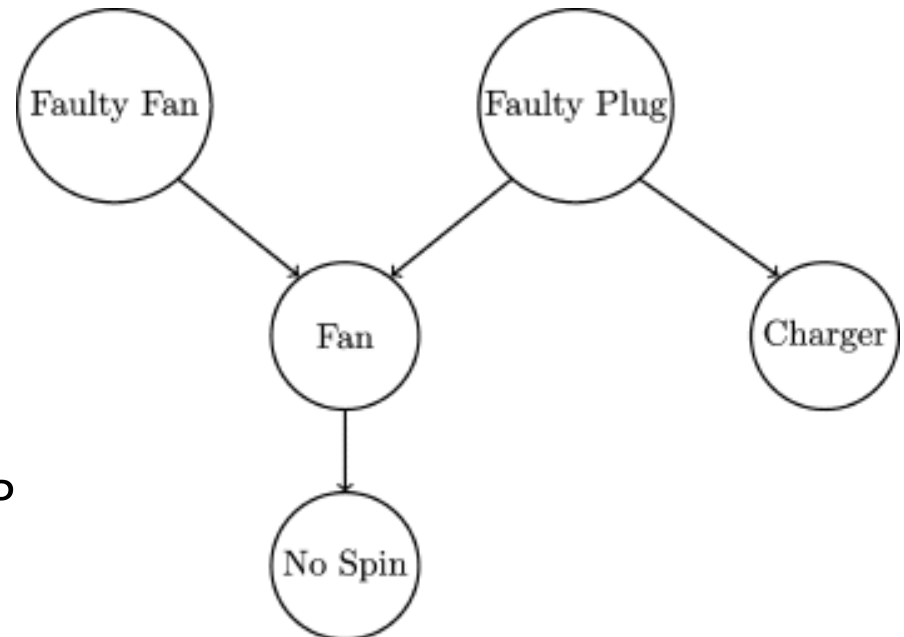
BN Example

Useful rule:

- Given the parents of a node A, the node A is independent of its non-descendants

Which are true?

- FF and FP are independent
- FF and FP independent given F
- F and C are independent
- F and C are independent given FP
- NS and C are independent
- NS and C are independent given FP
- NS and C are independent given F



Yes, no, no, yes, no, yes, yes

BN Example 2

Modelling the relationship between **Sprinkler**, **Rain** and **Grass Wet**

- Two events can cause the grass to become wet (**Rain=True** or **Sprinkler=True**)

BN Example 2

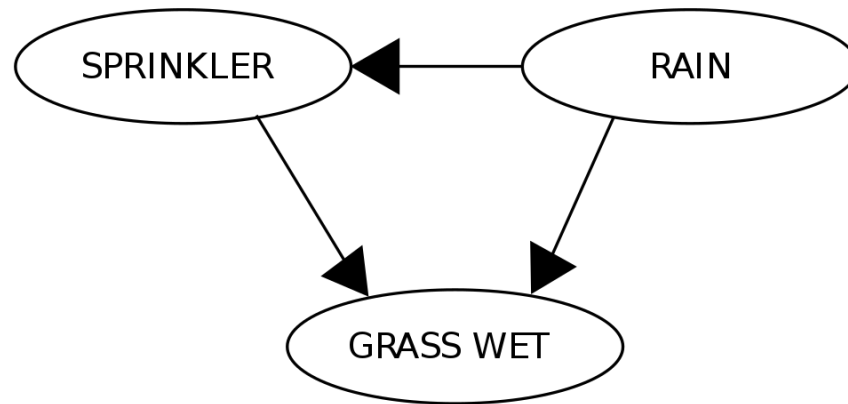
Modelling the relationship between **Sprinkler**, **Rain** and **Grass Wet**

- Two events can cause the grass to become wet (**Rain=True** or **Sprinkler=True**)
- If **Rain=True**, then Sprinkler is unlikely to be True

BN Example 2

Modelling the relationship between **Sprinkler**, **Rain** and **Grass Wet**

- Two events can cause the grass to become wet (**Rain=True** or **Sprinkler=True**)
- If **Rain=True**, then Sprinkler is unlikely to be True



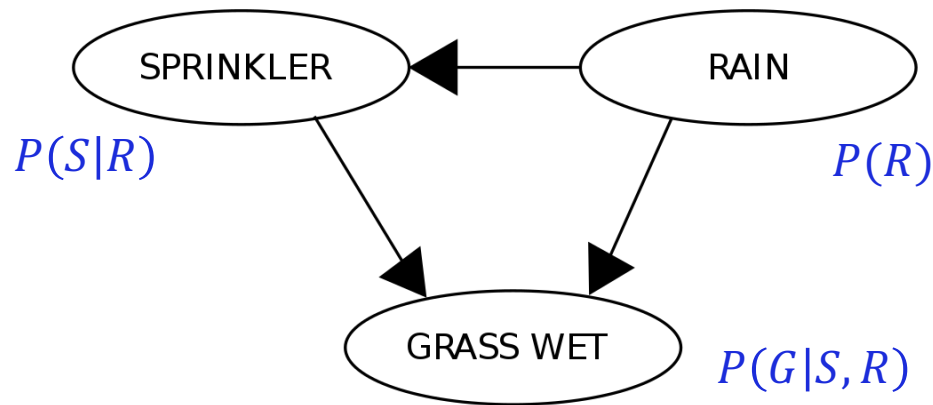
BN Example 2

Modelling the relationship between **Sprinkler**, **Rain** and **Grass Wet**

Let $G = \text{"Grass wet"}$, $S = \text{"Sprinkler turned on"}$, and $R = \text{"Raining"}$.

Joint probability is:

$$P(R, S, G) = P(G|S, R) P(S, R) = P(G|S, R)P(S|R)P(R)$$



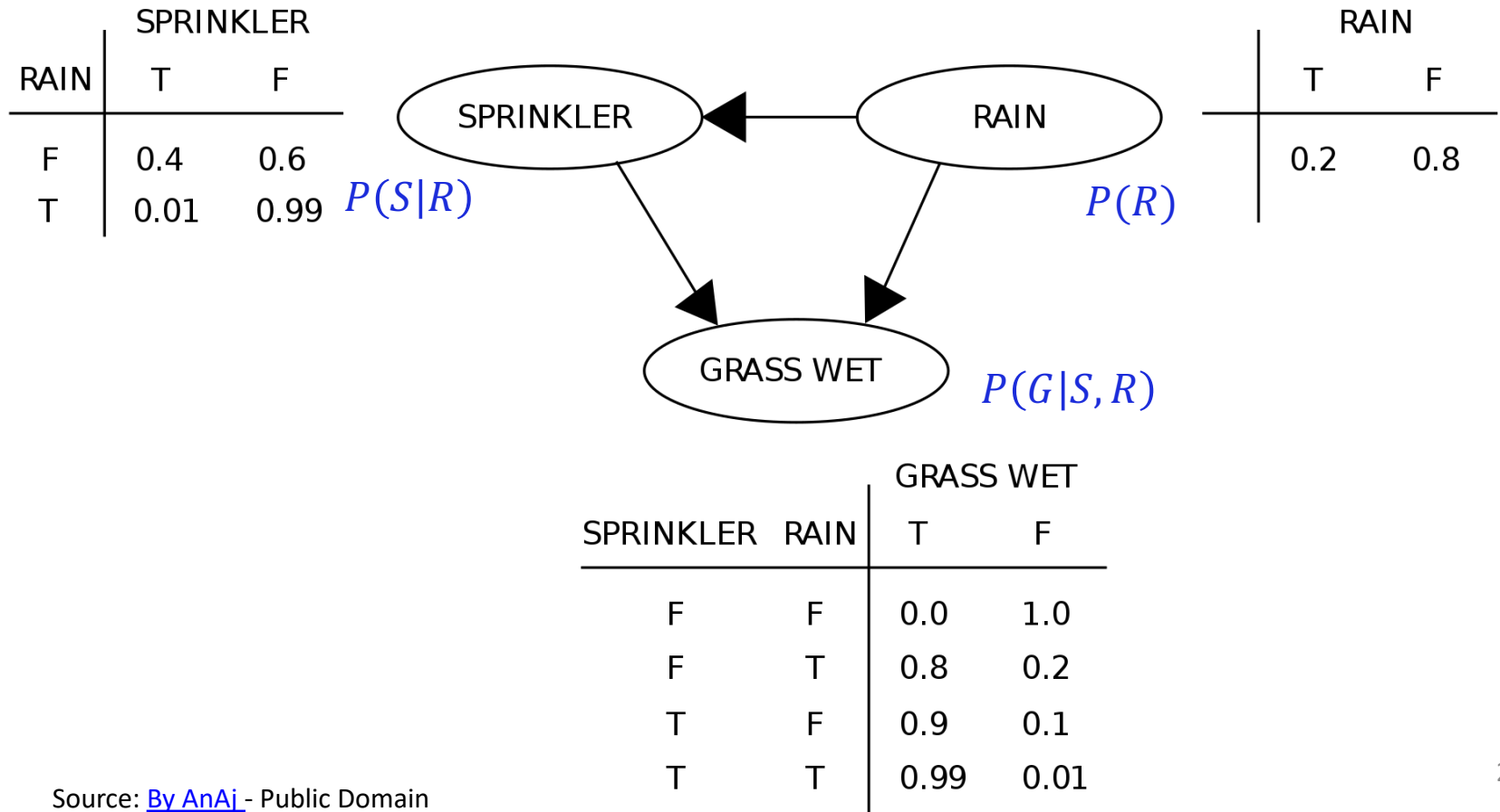
BN Example 2

Modelling the relationship between **Sprinkler**, **Rain** and **Grass Wet**

Let $G = \text{"Grass wet"}$, $S = \text{"Sprinkler turned on"}$, and $R = \text{"Raining"}$.

Joint probability is:

$$P(R, S, G) = P(G|S, R) P(S, R) = P(G|S, R)P(S|R)P(R)$$

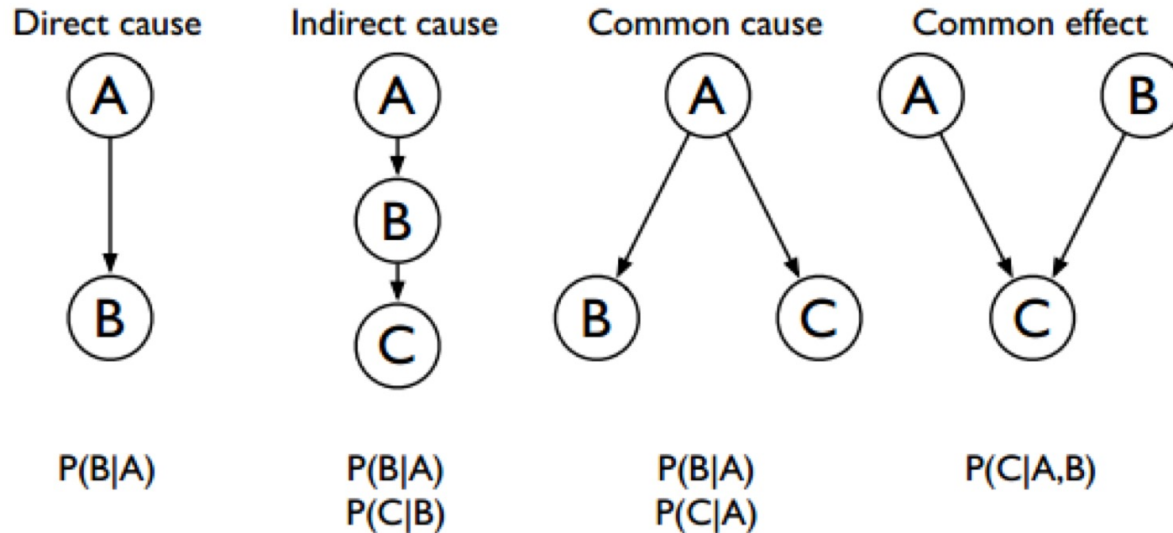


Why Bayesian Networks?

- The model encodes dependencies among all variables
- A BN can be used to learn causal relationships, and hence gain understanding about a problem domain and to predict the consequences of intervention
- The model has both a causal and probabilistic semantics, it is an ideal representation for combining prior knowledge (which often comes in causal form) and data
- Bayesian statistical methods in conjunction with Bayesian networks offer an efficient and principled approach for reduce the overfitting of data. (As it enforces a presumably correct structure.)
- A BN requires typically far fewer parameters and hence less storage than a full joint probability distribution

(In)Dependencies in BN

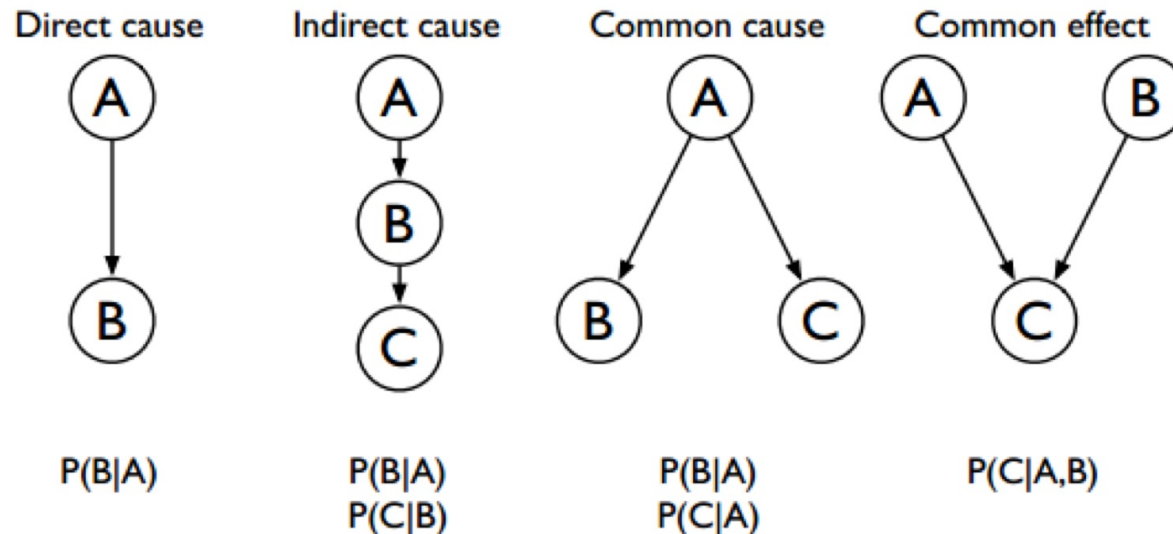
- **[Direct cause]:** A is a direct cause of B
 - A and B are **dependent**
- **[Indirect cause]:** A direct cause of B, B direct cause of C
 - A and C are **independent** given B



$$\begin{aligned}
 P(C|A) &= \sum_n P(C, B_n|A) \\
 &= \sum_n P(C|B_n, A) P(B_n|A) \\
 &= \sum_n P(C|B_n) P(B_n|A)
 \end{aligned}$$

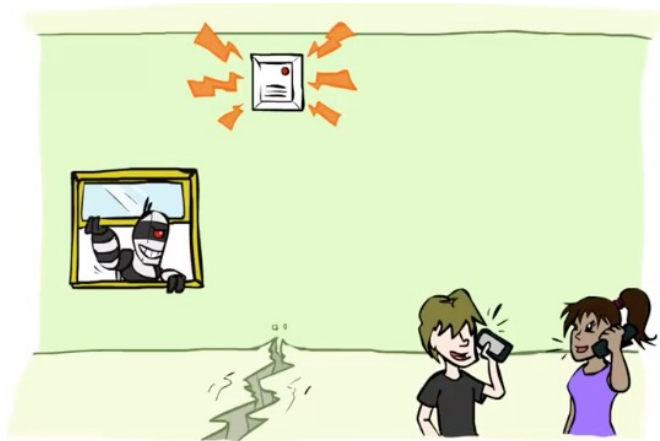
(In)Dependencies in BN

- **[Common Cause]:** A is a **common direct cause** of B and C
 - B and C are **dependent** (if A is not given)
 - B and C are **independent** given A
- **[Common Effect]:** C is a **common direct effect** of A and B
 - A and B are **independent** (if C is not given)
 - A and B are **dependent** given C (“explaining away”)



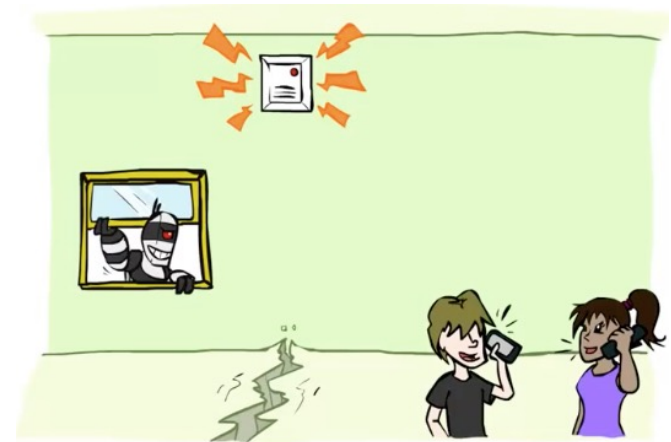
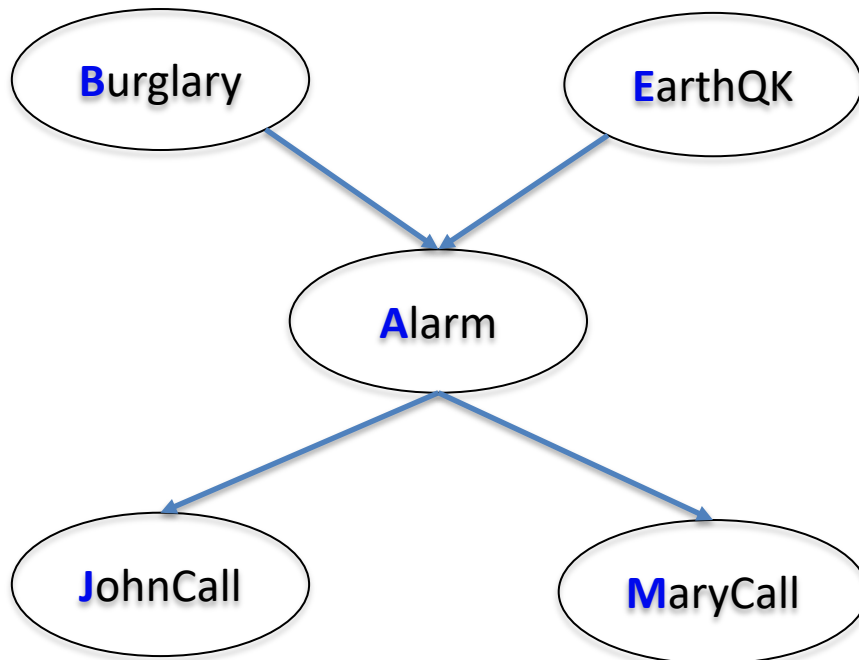
BN Example 3

- Your house has an **alarm** against **burglary**
 - The **alarm** will usually be set off by **burglars**
 - Sometimes it may also be set off by **earthquakes**
 - There are two neighbours, John and Mary
 - **John and Mary might call you** when they hear the alarm
 - They might also call you for other issues without alarm
- **Variables:**
 - Burglary, Earthquake, Alarm, JohnCalls, MaryCalls
 - All binary (**true** or **false**)
- **Relationship** between them?
 - Cause -> Effect



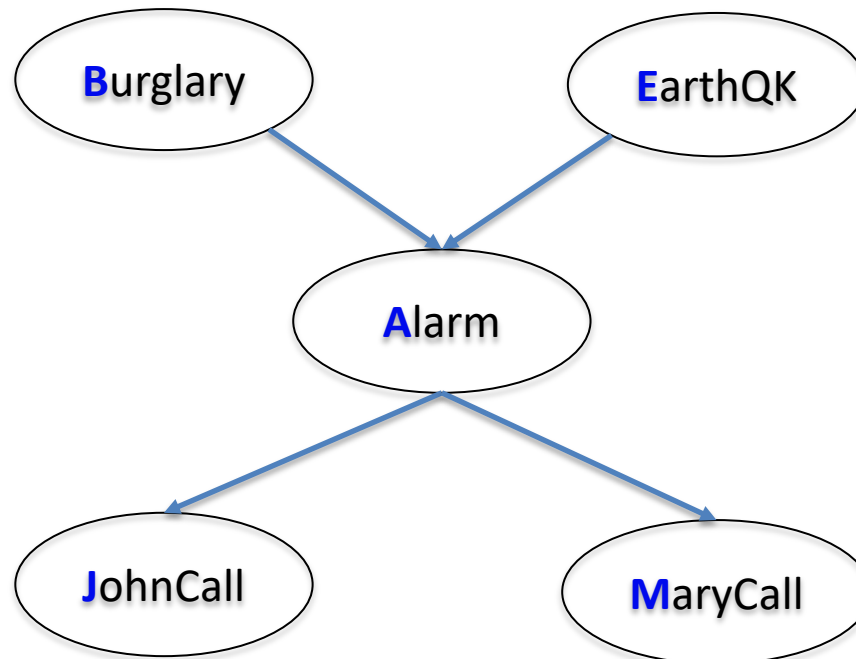
Alarm Network

- Domain **causal knowledge (causes and effects)**
 - A burglar can set the alarm off
 - An earthquake can set the alarm off
 - The alarm can cause Mary to call
 - The alarm can cause John to call



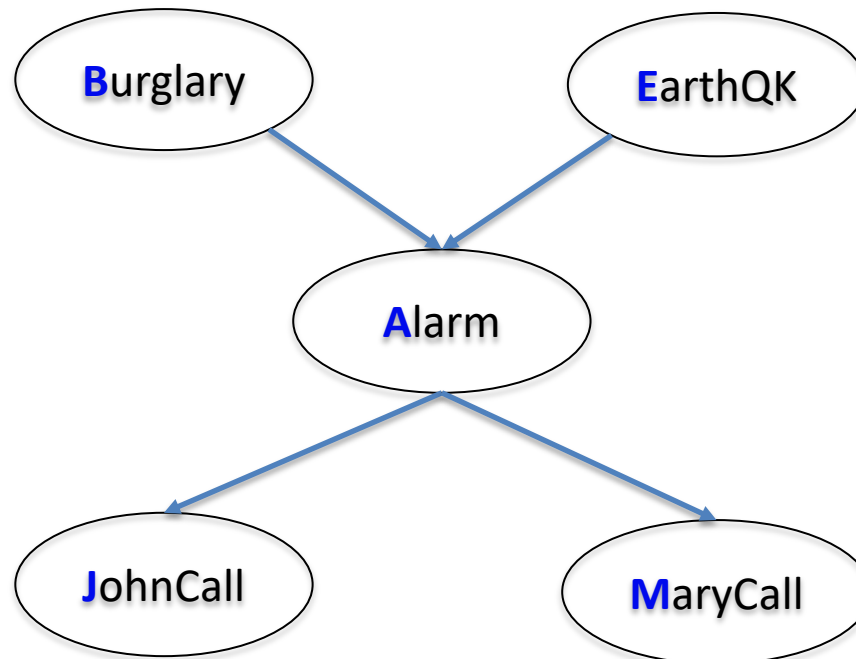
(In)Dependencies in BN

- Recall: given the parents of a node A, the node A is independent of its non-descendants



(In)Dependencies in BN

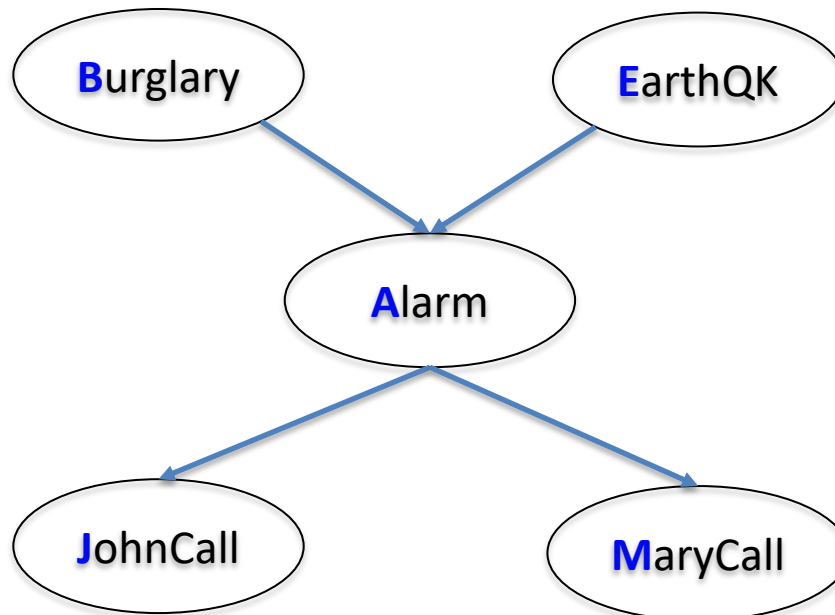
- Which are true?
 - B and E are independent
 - B and E are independent given A
 - B and M are independent
 - B and M are independent given A
 - J and M are independent
 - J and M are independent given A



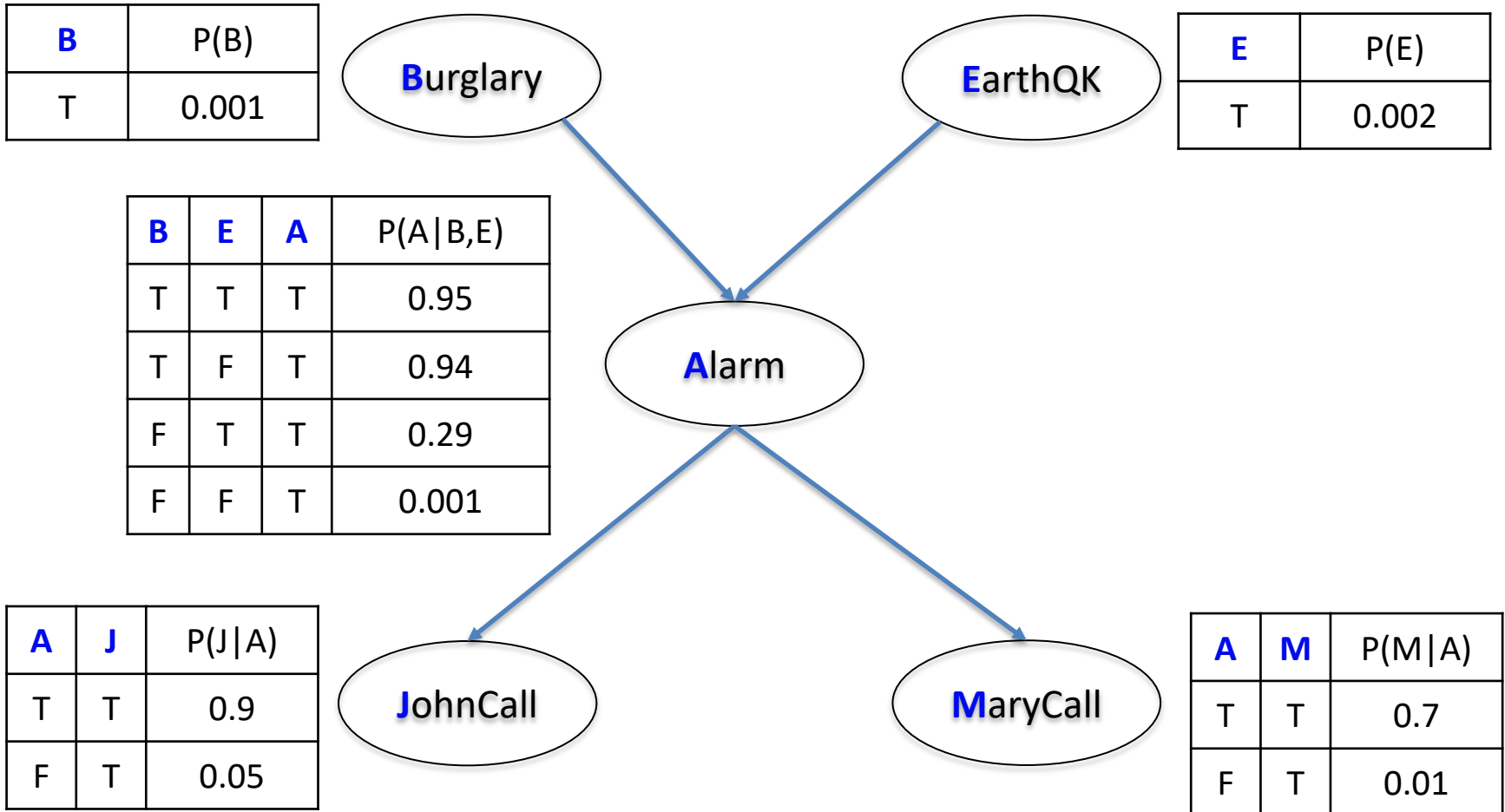
Factorisation

$$\begin{aligned}P(M, J, A, E, B) &= P(M, J, A | E, B)P(E, B) \\ &= P(M, J, A | E, B) P(E)P(B) \\ &= P(M, J | A, E, B)P(A | E, B)P(E)P(B) \\ &= P(M, J | A)P(A | E, B)P(E)P(B) \\ &= P(M | A)P(J | A)P(A | E, B)P(E)P(B)\end{aligned}$$

- Each conditional probability is a node.
- Table with 32 entries becomes tables with $2+2+4+1+1=10$ entries

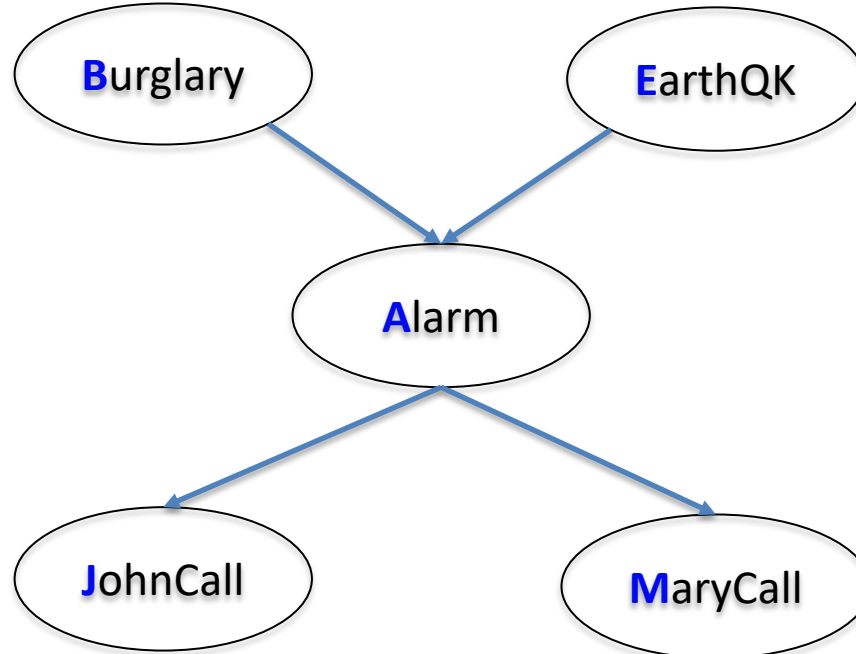


Alarm Network



Factorisation

- To factorise according to a Bayesian network: sort the variables so that the **causes are always before the effects**, e.g., [B, E, A, J, M], then use the rules:
 - $P(X_i | X_1, \dots, X_{i-1}) = P(X_i | \text{parents}(X_i), \dots, X_j, \dots)$
 $= P(X_i | \text{parents}(X_i))$
 - $P(X_i | X_j) = P(X_i)$ if X_i and X_j independent



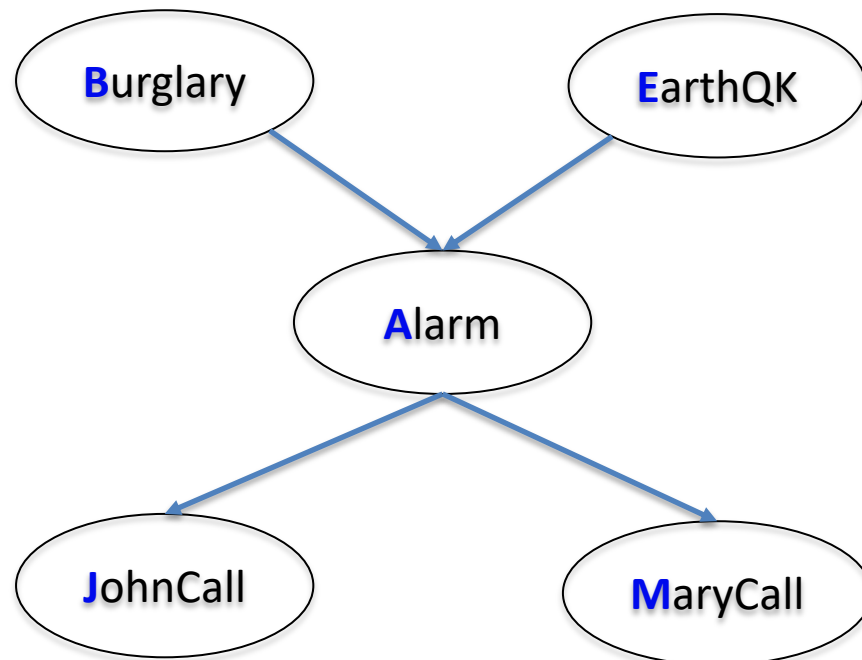
Factorisation

- We saw repeated application of the product rule gives

$$P(M, J, A, E, B) = P(M|A)P(J|A)P(A|E, B)P(E)P(B)$$

- The joint probability distribution over all variables in the network can be represented as a product of the conditional probabilities of each variable given its parents:

$$P(X_1, \dots, X_n) = P(X_n | \text{parents}(X_n))P(X_{n-1} | \text{parents}(X_{n-1})) \dots$$



Summary

- Conditionally independent given class label (NB)
- Bayes Net = Topology (graph) + Local Conditional Probabilities
- Factorisation

Coming up next...

- More on Bayesian Networks (Build a BN, # free parameters, ...)